

NetMHC-3.0: accurate web accessible predictions of human, mouse and monkey MHC class I affinities for peptides of length 8–11

Claus Lundegaard^{1,*}, Kasper Lamberth², Mikkel Harndahl², Søren Buus², Ole Lund¹ and Morten Nielsen¹

¹CBS, Department of Systems Biology, Technical University of Denmark DTU, Kemitorvet Build. 208, 2800 Lyngby

and ²Department of International Health, Immunology and Microbiology, University of Copenhagen, Panum Institute 22.3.6, Blegdamsvej 18, 2200 Copenhagen N, Denmark

Received January 31, 2008; Revised March 27, 2008; Accepted April 4, 2008

ABSTRACT

NetMHC-3.0 is trained on a large number of quantitative peptide data using both affinity data from the Immune Epitope Database and Analysis Resource (IEDB) and elution data from SYFPEITHI. The method generates high-accuracy predictions of major histocompatibility complex (MHC): peptide binding. The predictions are based on artificial neural networks trained on data from 55 MHC alleles (43 Human and 12 non-human), and position-specific scoring matrices (PSSMs) for additional 67 HLA alleles. As only the MHC class I prediction server is available, predictions are possible for peptides of length 8–11 for all 122 alleles. artificial neural network predictions are given as actual IC₅₀ values whereas PSSM predictions are given as a log-odds likelihood scores. The output is optionally available as download for easy post-processing. The training method underlying the server is the best available, and has been used to predict possible MHC-binding peptides in a series of pathogen viral proteomes including SARS, Influenza and HIV, resulting in an average of 75–80% confirmed MHC binders. Here, the performance is further validated and benchmarked using a large set of newly published affinity data, non-redundant to the training set. The server is free of use and available at: <http://www.cbs.dtu.dk/services/NetMHC>.

INTRODUCTION

Intracellular infections with pathogens such as viruses and certain bacteria are defeated by cytotoxic T lymphocytes (CTL). The CTL T-cell receptor (TCR) recognizes foreign

peptides in complex with major histocompatibility complex (MHC) class I molecules on the surface of the infected cells. MHC class I molecules preferably bind and present nine amino acid long peptides, which mainly originates from proteins expressed in the cytosol of the presenting cell. In most vertebrates, MHCs exist in a number of different allelic variants that each binds a specific and very limited set of peptides. For a number of years, prediction methods have developed to identify which peptides will bind a given MHC (1), and such predictions can be highly valuable in a broad range of applications, including rational vaccine design and disease diagnostics. The artificial neural network (ANN) training method behind NetMHC (2,3) has been benchmarked to be the best among available methods (4). Preliminary versions of the algorithm have been used to predict possible MHC-binding peptides in a large set of pathogenic viral proteomes, resulting in an average of >75% confirmed MHC binders (5). Most MHC prediction algorithms (a list of other servers is included in the Supplementary Material) are trained on peptides of the same length as they predict, but since data for peptide lengths different from nine are much more scarce, the broadness of MHC binding predictions for different peptide lengths is accordingly limited. In this server, however, a method is implemented making it possible to predict 8-, 10- and 11-mer peptide binding using 9-mer trained predictors, which extends the MHC coverage for these peptide lengths significantly compared to other available MHC:peptide-binding servers.

METHODS

The server is trained on the largest number of quantitative peptide:MHC affinity measurements ever published using both affinity data from the Immune Epitope Database and Analysis Resource (IEDB) (6), eluted peptide data from

*To whom correspondence should be addressed. Tel: +45 21900767; Fax: +45 45931585; Email: lunde@cbc.dtu.dk

the SYFPEITHI database (7) and proprietary affinity data. The predictions based on ANNs are trained essentially as described in (3) on data from 55 MHC alleles (43 Human and 12 non-human), and the predictions based on position specific scoring matrices (PSSMs) are trained as described in (2) for additional 67 HLA alleles. A large number of 9-mer MHC affinity data have become available from the IEDB database, since the training of the ANNs used at NetMHC-3.0, and all peptides not used in the training (6452 9-mer peptide affinity data points, covering 32 HLA alleles) were used for evaluation of the server performance. These data are available at the server. In this dataset, 3104 were measured to be binders ($IC_{50} < 500 \text{ nM}$), 76% of these were correctly predicted as such. 3030 peptides were predicted to bind to a given HLA, and 78% of these had a measured $IC_{50} < 500 \text{ nM}$. The average Pearson correlation coefficient (PCC) and area under a ROC curve (AUC) value using a 500 nM classification threshold were 0.71 and 0.86, respectively. For the full per allele results, see the Supplementary Material (Supplementary Table 1 and Supplementary Figure 1). NetMHC-3.0 uses a new approximation algorithm that reliably predicts the affinity of peptides of lengths 8, 10 and 11, for which affinity data for training are rare (8). The method uses predictors trained on peptides of length 9 to successfully extrapolate to other lengths. In short, the method approximates each peptide of any length to a number of 9-mers, by inserting X (for 8-mers) or deleting amino acid(s) (for 10- and 11-mers) and set the final prediction to an average of the 9-mer predictions. We had previously trained ANN predictors directly on 10-mer affinity data and since this training more than 2000 10-mer peptide:MHC affinities had become available from the IEDB database (6). Area under a ROC curve (AUC) values were calculated for each allele using either ANNs trained on 10-mers or the approximation method. For 12 of the 16 alleles, the approximation method performed better than the 10-mer trained ANNs ($P < 0.01$), see Supplementary Material Figure 2. However, for the four HLA-alleles, this evaluation showed better performance for ANNs trained on 10-mer peptides; these 10-mer trained ANNs are used for predictions by the server. For 8-mers, 2002 affinity data were extracted covering 35 MHC alleles. The overall PCC and AUC were 0.68 and 0.86, respectively. For 8-mer per allele performance, see the Supplementary Material Figure 4. For 8-mers, predictors trained on actual 8-mers seems to be better than the approximation method otherwise used, so for the alleles with available 8-mer affinity data, 8-mer trained ANNs are used for the predictions. In general, it is not possible to estimate how reliable a single prediction is. However, the stronger the affinity is predicted the higher are the chance that the actual affinity is stronger than the generally accepted binding threshold of 500 nM.

SERVER

NetMHC-3.0 predicts the binding affinity of either a list of peptides with a defined length (8–11 residues) or all possible sub-peptides hosted within full-length proteins.

```
>Protein_1
SLYNTVATLSLYNTATLSLYNTVATL
>Protein_2
ELEVENTENNINEANDEIGHTMERS
```

Figure 1. Example input in FASTA format.

The input must be in the FASTA format, or as peptides all of equal length, one peptide pr. line. The server will accept a maximum of 5000 sequences per submission; each sequence not more than 20 000 amino acids with a minimum length corresponding to the selected length of prediction (see subsequently). Input data can be pasted into a text field or uploaded from a local file on the user's computer.

If the input is in peptide format the corresponding tick-box must be selected. The input must not exceed 5000 sequences and with a maximum of 20 000 amino acids in each sequence. One or more MHCs must be selected, as well as the desired peptide length. Only one prediction length at a time can be used. The output can optionally be sorted according to the predicted affinity by selecting a tick-box. The predictions start by clicking the Submit button. An example input in FASTA format is shown in Figure 1.

The output is displayed as raw text with a header indicating the server name, the type of prediction (PSSM, ANN or ANN-approximation) the first selected allele and the date (Figure 2) followed by the prediction output in a column format. The columns are named in the first line of the prediction output. The first column [pos] is the position of the first amino acid of the predicted peptide within the possibly longer sequence, numbering starting with 0. Column (peptide) is the primary sequence of the (sub-)peptide. Column (logscore) is the raw prediction output, which for ANNs is $1 - \log_{50000}$ to the affinity in nanomolar units. For PSSM predictions the raw prediction score is a log-odds likelihood score. Additionally a column is included for ANN predictions, [affinity (nM)], which is the predicted affinity presented in nanomolar units. Column (Bind Level) indicates if the peptide is predicted to bind stronger than a certain threshold [for ANN predictions stronger than 50 nM (SB) or stronger than 500 nM (WB); for PSSM high-binding peptides (SB) have a prediction score greater than the 0.1% percentile score value of 1 000 000 random natural peptides, and weak binding (WB) peptides a score value above the 1% percentile score of 1 000 000 random natural peptides predictions]. Predicted affinities weaker than 500 nM or lower than the 1% percentile score have no indications. Column (Protein Name) gives the name of the predicted protein. If peptide input was used, the name will always be 'Sequence'. Column (Allele) gives the name of the MHC allele chosen. The output contains all the sub-peptides for each protein for a given allele either in the order they appear in the sequence or sorted by predicted affinity within each protein (if chosen). If more than one protein sequence were entered, a dashed line will separate the

NetMHC 3.0 Server – prediction results

http://www.cbs.dtu.dk/cgi-bin/nph-webface?jobid=netmhc_47A073EE0A46069F

Keystone Influenza virus resource Allelefrequencies NIH A-C list Yahoo! CBS Python dict types Nucleic Ac...ERVER ISSUE

NetMHC 3.0 Server - prediction results

Technical University of Denmark

Wednesday January 30 2008 13:56

NetMHC version 3.0. 10mer predictions using Artificial Neural Networks approximation. Allele A0201. Strong binder threshold 50 nM. Weak binder threshold score 500 nM

[Download output sheet](#)

pos	peptide	logscore	affinity(nM)	Bind Level	Protein Name	Allele
9	SLYNTATLSL	0.677	33	SB	Protein_1	A0201
7	TLSLYNTATL	0.527	166	WB	Protein_1	A0201
15	TLSLYNTVAT	0.314	1673		Protein_1	A0201
16	LSLYNTVATL	0.270	2685		Protein_1	A0201
14	ATLSLYNTVA	0.242	3649		Protein_1	A0201
0	SLYNTVATLS	0.212	5033		Protein_1	A0201
13	TATLSLYNTV	0.206	5372		Protein_1	A0201
4	TVATLSLYNT	0.166	8311		Protein_1	A0201
6	ATLSLYNTAT	0.145	10400		Protein_1	A0201
12	NTATLSLYNT	0.106	15909		Protein_1	A0201
1	LYNTVATLSL	0.105	15994		Protein_1	A0201
5	VATLSLYNTA	0.087	19487		Protein_1	A0201
8	LSLYNTATLS	0.047	30200		Protein_1	A0201
3	NTVATLSLYN	0.024	38693		Protein_1	A0201
11	YNTATLSLYN	0.021	39650		Protein_1	A0201
2	YNTVATLSLY	0.009	45508		Protein_1	A0201
10	LYNTATLSLY	0.007	46459		Protein_1	A0201
12	EANDEIGHTM	0.029	36689		Protein_2	A0201
9	NINEANDEIG	0.024	38650		Protein_2	A0201
1	LEVENTENNI	0.023	39109		Protein_2	A0201
8	NNINEANDEI	0.023	39177		Protein_2	A0201
4	ENTENNINNEA	0.008	45947		Protein_2	A0201
11	NEANDEIGHT	0.008	45951		Protein_2	A0201
13	ANDEIGHTME	0.008	46055		Protein_2	A0201
0	ELEVENTENN	0.007	46354		Protein_2	A0201
5	NTENNINNEAN	0.006	46901		Protein_2	A0201
6	TENNINNEAND	0.004	47719		Protein_2	A0201
3	VENTENNINNE	0.004	47728		Protein_2	A0201
15	DEIGHTMERS	0.002	48666		Protein_2	A0201
10	INEANDEIGH	0.002	49001		Protein_2	A0201
14	NDEIGHTMER	0.001	49293		Protein_2	A0201
2	EVENTENNIN	0.001	49313		Protein_2	A0201
7	ENNINNEANDE	0.001	49640		Protein_2	A0201
NetMHC version 3.0. 10mer predictions using Artificial Neural Networks. Allele A0301.						
Strong binder threshold 50 nM. Weak binder threshold score 500 nM						
Download output sheet						
pos	peptide	logscore	affinity(nM)	Bind Level	Protein Name	Allele
10	LYNTATLSLY	0.450	382	WB	Protein_1	A0301
9	SLYNTATLSL	0.284	2312		Protein_1	A0301
0	SLYNTVATLS	0.278	2477		Protein_1	A0301
2	YNTVATLSLY	0.240	3719		Protein_1	A0301
8	LSLYNTATLS	0.140	11049		Protein_1	A0301
3	NTVATLSLYN	0.135	11580		Protein_1	A0301
15	TLSLYNTVAT	0.109	15383		Protein_1	A0301

Figure 2. Raw text output using the input in Figure 1 and selecting the alleles HLA-A0201 and HLA-A0301. 10-mer peptide predictions were chosen. Affinity sorting was chosen.

The screenshot shows a Microsoft Excel spreadsheet titled "tmpQSgZP5.xls". The columns are labeled A through G. Column A contains protein names like "Protein_1" and "Protein_2". Column B contains peptide sequences. Column C contains positions. Column D contains "A0201 ANN Approximation predicted affinity (Kd, nM)". Column E contains "A0301 ANN Direct predicted affinity (Kd, nM)". Column F contains "Average score (higher score = stronger affinity)". The data is organized by protein and peptide sequence.

Figure 3. Downloaded output sheet opened in Microsoft® Excel and adjusted with of column. The output was generated using input in Figure 1 and selecting the alleles HLA-A0201 and HLA-A0301. 10mer peptide predictions were chosen.

peptides from each protein. If more than one allele were chosen, the output will show a header similar to the first immediately after the first predictions, all in the same web output page.

In each header, there is a link to a file with the output in tab as separated format, where the filename ends on.xls making it easily imported into spreadsheet programs. This file always contains the predicted peptides in the order they appeared in the input file. The output data for each peptide will be displayed on a single line with predictions for each of the selected alleles in different columns (Figure 3).

FINAL REMARKS

This server is developed to aid research and limit the resources needed for rational and effective CTL epitope discovery and will be continuously updated as new data become available. All comments and suggestions for usability improvements are most welcome.

SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

ACKNOWLEDGEMENTS

This work was funded by European Commission (LSHB-CT-2003-503231, LSHB-CT-2004-012175) and National

Institutes of Health (HHSN26600400006C, HHSN26600400025C, HHSN266200400083C).

Conflict of interest statement. None declared.

REFERENCES

- Lundsgaard,C., Lund,O., Kesmir,C., Brunak,S. and Nielsen,M. (2007) Modeling the adaptive immune system: predictions and simulations. *Bioinformatics*, **23**, 3265–3275.
- Nielsen,M., Lundsgaard,C., Worning,P., Hvid,C.S., Lamberth,K., Buus,S., Brunak,S. and Lund,O. (2004) Improved prediction of MHC class I and class II epitopes using a novel Gibbs sampling approach. *Bioinformatics*, **20**, 1388–1397.
- Nielsen,M., Lundsgaard,C., Worning,P., Lauemoller,S.L., Lamberth,K., Buus,S., Brunak,S. and Lund,O. (2003) Reliable prediction of T-cell epitopes using neural networks with novel sequence representations. *Proteins Sci*, **12**, 1007–1017.
- Peters,B., Bui,H.H., Frankild,S., Nielsen,M., Lundsgaard,C., Kostem,E., Basch,D., Lamberth,K., Harndahl,M., Flerl,W. et al. (2006) A community resource benchmarking predictions of peptide binding to MHC-I molecules. *PLoS Comput Biol*, **2**, e65.
- Sylvester-Hvid,C., Nielsen,M., Lamberth,K., Roder,G., Justesen,S., Lundsgaard,C., Worning,P., Thomadsen,H., Lund,O., Brunak,S. et al. (2004) SARS CTL vaccine candidates; HLA supertype-, genome-wide scanning and biochemical validation. *Tissue Antigens*, **63**, 395–400.
- Sette,A., Flerl,W., Peters,B., Sathiamurthy,M., Bui,H.H. and Wilson,S. (2005) A roadmap for the immunomics of category A-C pathogens. *Immunity*, **22**, 155–161.
- Rammensee,H., Bachmann,J., Emmerich,N.P., Bachor,O.A. and Stevanovic,S. (1999) SYFPEITHI: database for MHC ligands and peptide motifs. *Immunogenetics*, **50**, 213–219.
- Lundsgaard,C., Lund,O. and Nielsen,M. (2008) Accurate approximation method for prediction of class I MHC af-finities for peptides of length 8, 10 and 11 using prediction tools trained on 9mers. *Bioinformatics*, in press, doi:10.1093/bioinformatics/btn128.